# Relational Database Design Theory

Introduction to Databases CompSci 316 Spring 2017

DUKE COMPUTER SCIENCE



• Homework #1 due Monday 02/06 (11:59 pm)



 Wouldn't it be nice to have a systematic approach to detecting and removing redundancy in designs?







### Redefining "keys" using FD's

- A set of attributes K is a key for a relation R if
- K → all (other) attributes of R
  That is, K is a "super key"
- No proper subset of K satisfies the above condition
  That is, K is minimal

### Reasoning with FD's

Given a relation R and a set of FD's  $\mathcal{F}$ 

- Does another FD follow from  $\mathcal{F}?$ 
  - Are some of the FD's in  ${\mathcal F}$  redundant (i.e., they follow from the others)?
- Is K a key of R?
  What are all the keys of R?

#### Attribute closure

- Given R, a set of FD's  $\mathcal{F}$  that hold in R, and a set of attributes Z in R: The closure of Z (denoted  $Z^+$ ) with respect to  $\mathcal{F}$  is the set of all attributes  $\{A_1, A_2, ...\}$  functionally determined by Z (that is,  $Z \to A_1A_2$  ...)
- Algorithm for computing the closure
  - Start with closure = Z
  - If  $X \to Y$  is in  $\mathcal{F}$  and X is already in the closure, then also add Y to the closure
  - Repeat until no new attributes can be added

#### A more complex example

UserJoinsGroup (uid, uname, twitterid, gid, fromDate) Assume that there is a 1-1 correspondence between our users and Twitter accounts

- uid  $\rightarrow$  uname, twitterid
- twitterid  $\rightarrow$  uid
- uid, gid  $\rightarrow$  fromDate

Not a good design, and we will see why shortly





### Rules of FD's

- Armstrong's axioms
  - Reflexivity: If  $Y \subseteq X$ , then  $X \to Y$
  - Augmentation: If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$  for any Z
  - Transitivity: If  $X \to Y$  and  $Y \to Z$ , then  $X \to Z$
- Rules derived from axioms
  - Splitting: If  $X \to YZ$ , then  $X \to Y$  and  $X \to Z$
  - Combining: If  $X \to Y$  and  $X \to Z$ , then  $X \to YZ$
- ${}^{\mathscr{F}}$  Using these rules, you can prove or disprove an FD given a set of FDs



<sup>®</sup> Example of redundancy							15
UserJoinsGroup (uid, uname, twitterid, gid, fromDate) • uid → uname, twitterid ( plus other FD's)							
	uid	uname	twitterid	gid	fromDate		
	142	Bart	@BartJSimpson	dps	1987-04-19		
	123	Milhouse	@MilhouseVan_	gov	1989-12-17		
	857	Lisa	@lisasimpson	abc	1987-04-19		
	857	Lisa	@lisasimpson	gov	1988-09-01		
	456	Ralph	@ralphwiggum	abc	1991-04-25		
	456	Ralph	@ralphwiggum	gov	1992-09-01		
How can we eliminate the redundancy?							y?



Unnecessary decomposition	1;



## Lossless join decomposition

- Decompose relation R into relations S and T
  - $attrs(R) = attrs(S) \cup attrs(T)$
  - $S = \pi_{attrs(S)}(R)$
  - $T = \pi_{attrs(T)}(R)$
- The decomposition is a lossless join decomposition if, given known constraints such as FD's, we can guarantee that  $R = S \bowtie T$



### Questions about decomposition

- When to decompose
- How to come up with a correct decomposition (i.e., lossless join decomposition)

#### An answer: BCNF

- A relation *R* is in Boyce-Codd Normal Form if
   For every non-trivial FD *X* → *Y* in *R*, *X* is a super key
  - That is, all FDs follow from "key  $\rightarrow$  other attributes"

When to decompose
As long as some relation is not in BCNF

- How to come up with a correct decomposition
  - Always decompose on a BCNF violation (details next)
     Then it is guaranteed to be a lossless join decomposition!

### BCNF decomposition algorithm

• Find a BCNF violation

- That is, a non-trivial FD  $X \to Y$  in R where X is not a super key of R
- Decompose *R* into *R*<sub>1</sub> and *R*<sub>2</sub>, where
  - $R_1$  has attributes  $X \cup Y$
  - $R_2$  has attributes  $X \cup Z$ , where Z contains all attributes of R that are in neither X nor Y
- Repeat until all relations are in BCNF

## BCNF decomposition example

uid  $\rightarrow$  uname, twitterid twitterid  $\rightarrow$  uid uid, gid  $\rightarrow$  fromDate

UserJoinsGroup (uid, uname, twitterid, gid, fromDate)

<sup>Any decomposition gives R ⊆ S ⋈ T (why?)
A lossy decomposition is one with R ⊂ S ⋈ T</sup> 

#### Another example

uid  $\rightarrow$  uname, twitterid twitterid  $\rightarrow$  uid uid, gid  $\rightarrow$  fromDate

UserJoinsGroup (uid, uname, twitterid, gid, fromDate)

#### Why is BCNF decomposition lossless

Given non-trivial  $X \rightarrow Y$  in R where X is not a super key of R, need to prove:

- Anything we project always comes back in the join:  $R \subseteq \pi_{XY}(R) \bowtie \pi_{XZ}(R)$ 
  - Sure; and it doesn't depend on the FD
- Anything that comes back in the join must be in the original relation:
  - $R \supseteq \pi_{XY}(R) \bowtie \pi_{XZ}(R)$
  - Proof will make use of the fact that  $X \to Y$

#### Recap

- Functional dependencies: a generalization of the key concept
- Non-key functional dependencies: a source of redundancy
- BCNF decomposition: a method for removing redundancies
  - BNCF decomposition is a lossless join decomposition
- BCNF: schema in this normal form has no redundancy due to FD's

#### BCNF = no redundancy?

- User (uid, gid, place)
  - A user can belong to multiple groups
  - A user can register places she's visited
  - Groups and places have nothing to do with other

uid	gid	place
142	dps	Springfield
142	dps	Australia
456	abc	Springfield
456	abc	Morocco
456	gov	Springfield
456	gov	Morocco

#### Multivalued dependencies

- A multivalued dependency (MVD) has the form *X* → *Y*, where *X* and *Y* are sets of attributes in a relation *R*
- X → Y means that whenever two rows in R agree on all the attributes of X, then we can swap their Y components and get two rows that are also in R



### **MVD** examples

User (uid, gid, place)

- uid → gid
- uid → place
- Intuition: given uid, gid and place are "independent"
- uid, gid  $\rightarrow$  place
  - Trivial: LHS U RHS = all attributes of R
- uid, gid → uid
  Trivial: LHS ⊇ RHS

#### Complete MVD + FD rules

- FD reflexivity, augmentation, and transitivity
- MVD complementation:
- If  $X \rightarrow Y$ , then  $X \rightarrow attrs(R) X Y$
- MVD augmentation: If  $X \twoheadrightarrow Y$  and  $V \subseteq W$ , then  $XW \twoheadrightarrow YV$
- MVD transitivity: If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z - Y$
- Replication (FD is MVD): If  $X \to Y$ , then  $X \to Y$ Try proving things using these!?
- Coalescence: If  $X \rightarrow Y$  and  $Z \subseteq Y$  and there is some W disjoint from Y such that  $W \rightarrow Z$ , then  $X \rightarrow Z$

### An elegant solution: chase

- Given a set of FD's and MVD's  $\mathcal{D}$ , does another dependency d (FD or MVD) follow from  $\mathcal{D}$ ?
- Procedure
  - Start with the premise of *d*, and treat them as "seed" tuples in a relation
  - Apply the given dependencies in  $\mathcal D$  repeatedly If we apply an FD, we infer equality of two symbols • If we apply an MVD, we infer more tuples
  - If we infer the conclusion of *d*, we have a proof
  - Otherwise, if nothing more can be inferred, we have a counterexample









### 4NF decomposition algorithm

- Find a 4NF violation
  - A non-trivial MVD  $X \rightarrow Y$  in R where X is not a superkey
- Decompose *R* into *R*<sub>1</sub> and *R*<sub>2</sub>, where
  - $R_1$  has attributes  $X \cup Y$
  - $R_2$  has attributes  $X \cup Z$  (where Z contains R attributes not in X or Y)
- Repeat until all relations are in 4NF
- Almost identical to BCNF decomposition algorithm
- Any decomposition on a 4NF violation is lossless

## 4NF decomposition example

uid

#### Summary

- Philosophy behind BCNF, 4NF: Data should depend on the key, the whole key, and nothing but the key!
  - You could have multiple keys though
- Other normal forms
  - 3NF: More relaxed than BCNF; will not remove redundancy if doing so makes FDs harder to enforce
  - 2NF: Slightly more relaxed than 3NF
  - 1NF: All column values must be atomic